

Chapter 6

TIME-MARCHING METHODS FOR ODE'S

In the following we use u as for the general representation of the dependent variable. If we wish to emphasize that our analysis is being carried out in eigenspace, then w will be used instead. Using the semi-discrete approach, we reduce our PDE to a set of coupled ODE's represented in general by Eq. 4.3. However, for the purpose of this chapter, we need only consider the scalar case

$$\frac{du}{dt} = u' = F(u, t) \quad (6.1)$$

Our first task is to find numerical approximations that can be used to carry out the time integration of Eq. 6.1 to some given accuracy, where accuracy can be measured either in a local or a global sense. We then face a further task concerning the numerical stability of the resulting methods, but we postpone such considerations to the next chapter.

In Chapter 2, we introduced the convention that the n subscript, or the (n) superscript, always points to a discrete time value, and h represents the time interval Δt . Combining this notation with Eq. 6.1 gives

$$u'_n = F_n = F(u_n, t_n) \quad ; \quad t = nh$$

Often we need a more sophisticated notation for intermediate time steps involving temporary calculations denoted by \tilde{u} , \bar{u} , etc. For these we use the notation

$$\tilde{u}'_{n+\alpha} = \tilde{F}_{n+\alpha} = F(\tilde{u}_{n+\alpha}, t_n + \alpha h)$$

The choice of u' or F to express the derivative in a scheme is arbitrary. They are both commonly used in the literature on ODE's.

The methods we study are to be applied to linear or nonlinear ODE's, but the methods themselves are formed by linear combinations of the dependent variable and its derivative at various time intervals. They are represented conceptually by

$$u_{n+1} = f(\beta_1 u'_{n+1}, \beta_0 u'_n, \beta_{n-1} u'_{n-1}, \dots, \alpha_0 u_n, \alpha_{-1} u_{n-1}, \dots) \quad (6.2)$$

With an appropriate choice of the α 's and β 's, these methods can be constructed to give a local Taylor series accuracy of any order. The methods are said to be *explicit* if $\beta_1 = 0$ and *implicit* otherwise. An *explicit* method is one in which the new predicted solution is only a function of known data, for example, u'_n, u'_{n-1}, u_n , and u_{n-1} for a method using two previous time levels, and therefore the time advance is simple. For an *implicit* method, the new predicted solution is also a function of the time derivative at the new time level, that is, u'_{n+1} . As we shall see, for systems of ODE's and nonlinear problems, *implicit* methods require more complicated strategies to solve for u_{n+1} than *explicit* methods.

6.1 Converting Time-Marching Methods to O Δ E 's

Examples of some very common forms of methods used for time-marching general ODE's are:

$$u_{n+1} = u_n + h u'_n \quad (6.3)$$

$$u_{n+1} = u_n + h u'_{n+1} \quad (6.4)$$

and

$$\begin{aligned} \tilde{u}_{n+1} &= u_n + h u'_n \\ u_{n+1} &= \frac{1}{2} [u_n + \tilde{u}_{n+1} + h \tilde{u}'_{n+1}] \end{aligned} \quad (6.5)$$

According to the conditions presented under Eq. 6.2, the first and third of these are examples of explicit methods. We refer to them as the Euler method and the MacCormack predictor-corrector method, respectively. The second is implicit and referred to as the implicit (or backward) Euler method.

These methods are simple recipes for the time advance of a function in terms of its value and the value of its derivative, at given time intervals. The material presented in Chapter 4 develops a way to evaluate such methods by introducing the concept of the representative equation

$$\frac{du}{dt} = \lambda u + a e^{\mu t} \quad (6.6)$$

written here in terms of the dependent variable, u . The value of this equation arises from the fact that it permits us to convert a linear time-marching method into a linear ordinary *difference* equation (OΔE). The latter are subject to a whole body of analysis that is similar in many respects to, and just as powerful as, the theory of ordinary differential equations, (ODE's). We next consider examples of this conversion process and then go into the general theory on solving OΔE's.

Apply the simple explicit Euler scheme, Eq. 6.3, to Eq. 6.6. There results

$$u_{n+1} = u_n + h(\lambda u_n + ae^{\mu hn})$$

or

$$u_{n+1} - (1 + \lambda h)u_n = hae^{\mu hn} \quad (6.7)$$

Eq. 6.7 is a linear ordinary OΔE with constant coefficients expressed in terms of the dependent variable u_n and the independent variable n . As another example, applying the implicit Euler method, Eq. 6.4, to Eq. 6.6, we find

$$u_{n+1} = u_n + h(\lambda u_{n+1} + e^{\mu h(n+1)})$$

or

$$(1 - \lambda h)u_{n+1} - u_n = he^{\mu h} \cdot ae^{\mu hn} \quad (6.8)$$

As a final example, the predictor-corrector sequence given in Eq. 6.5 gives

$$\begin{aligned} \tilde{u}_{n+1} - (1 + \lambda h)u_n &= ahe^{\mu hn} \\ -\frac{1}{2}(1 + \lambda h)\tilde{u}_{n+1} + u_{n+1} - \frac{1}{2}u_n &= \frac{1}{2}ahe^{\mu h(n+1)} \end{aligned} \quad (6.9)$$

which is a coupled set of linear OΔE's with constant coefficients.

Now we need to develop techniques for analyzing these difference equations so that we can compare the merits of the time-marching methods that generated them.

6.2 Solution of Linear OΔE's With Constant Coefficients

The techniques for solving *linear difference equations* with constant coefficients is as well developed as that for ODE's and the theory follows a remarkably parallel path. This is demonstrated by repeating some of the developments in Section 4.3, but for difference rather than differential equations.

6.2.1 First- and Second-Order Difference Equations

First-Order Equations

The simplest nonhomogeneous OΔE of interest is given by the single first-order equation

$$u_{n+1} = \sigma u_n + ab^n \quad (6.10)$$

where σ , a , and b are, in general, complex parameters. The independent variable is now n rather than t , and since the equations are linear and have constant coefficients (i.e., are essentially autonomous), σ is not a function of either n or u . The exact solution of Eq. 6.10 is

$$u_n = c_1(\sigma)^n + \frac{ab^n}{b - \sigma}$$

where c_1 is a constant determined by the initial conditions. In terms of the initial value of u it can be written

$$u_n = u_0\sigma^n + a\frac{b^n - \sigma^n}{b - \sigma}$$

Just as in the development of Eq. 4.10, one can readily show that the solution of the defective case, ($b = \sigma$),

$$u_{n+1} = \sigma u_n + a\sigma^n$$

is

$$u_n = \left[u_0 + an\sigma^{-1} \right] \sigma^n$$

This can all be easily verified by substitution.

Second-Order Equations

The homogeneous form of a second-order difference equation is given by

$$u_{n+2} + a_1u_{n+1} + a_0u_n = 0 \quad (6.11)$$

Instead of the differential operator $D \equiv \frac{d}{dt}$ used for ODE's, we use for OΔE's the difference operator E (commonly referred to as the displacement or shift operator) and defined formally by the relations

$$u_{n+1} = Eu_n \quad , \quad u_{n+k} = E^k u_n$$

. Further notice that the displacement operator also applies to exponents, thus

$$b^\alpha \cdot b^n = b^{n+\alpha} = E^\alpha \cdot b^n$$

where α can be any fraction or irrational number.

The roles of D and E are the same insofar as once they have been introduced to the basic equations the value of $u(t)$ or u_n can be factored out. Thus Eq. 6.11 can now be re-expressed in an operational notion as

$$(E^2 + a_1 E + a_0)u_n = 0 \quad (6.12)$$

which must be zero for all u_n . Eq. 6.12 is known as the *operational form* of Eq. 6.11. The operational form contains a characteristic polynomial $P(E)$ which plays the same role for difference equations that $P(D)$ played for differential equations; that is, its roots determine the solution to the OΔE. In the analysis of OΔE's, we label these roots $\sigma_1, \sigma_2, \dots$, etc, and refer to them as the σ -roots. They are found by solving the equation $P(\sigma) = 0$. In the simple example given above, there are just two σ roots and in terms of them the solution can be written

$$u_n = c_1(\sigma_1)^n + c_2(\sigma_2)^n \quad (6.13)$$

where c_1 and c_2 depend upon the initial conditions. The fact that Eq. 6.13 is a solution to Eq. 6.11 for all c_1, c_2 and n should be verified by substitution.

6.2.2 Special Cases of Coupled First-Order Equations

A Complete System

Coupled, first-order, linear homogeneous difference equations have the form

$$\begin{aligned} u_1^{(n+1)} &= c_{11}u_1^{(n)} + c_{12}u_2^{(n)} \\ u_2^{(n+1)} &= c_{21}u_1^{(n)} + c_{22}u_2^{(n)} \end{aligned} \quad (6.14)$$

which can also be written

$$\vec{u}_{n+1} = C\vec{u}_n \quad , \quad \vec{u}_n = [u_1^{(n)}, u_2^{(n)}]^T \quad , \quad C = \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

The operational form of Eq. 6.14 can be written

$$\begin{bmatrix} (c_{11} - E) & c_{12} \\ c_{21} & (c_{22} - E) \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}^{(n)} = [C - E I]\vec{u}_n = 0$$

which must be zero for all u_1 and u_2 . Again we are led to a characteristic polynomial, this time having the form $P(E) = \det [C - E I]$. The σ -roots are found from

$$P(\sigma) = \det \begin{bmatrix} (c_{11} - \sigma) & c_{12} \\ c_{21} & (c_{22} - \sigma) \end{bmatrix} = 0$$

Obviously the σ_k are the eigenvalues of C and, following the logic of Section 4.3, if \vec{x} are its eigenvectors, the solution of Eq. 6.14 is

$$\vec{u}_n = \sum_{k=1}^2 c_k (\sigma_k)^n \vec{x}_k$$

where c_k are constants determined by the initial conditions.

A Defective System

The solution of O Δ E's with defective eigensystems follows closely the logic in Section 4.3.2 for defective ODE's. For example, one can show that the solution to

$$\begin{bmatrix} \bar{u}_{n+1} \\ \hat{u}_{n+1} \\ u_{n+1} \end{bmatrix} = \begin{bmatrix} \sigma & & \\ 1 & \sigma & \\ & 1 & \sigma \end{bmatrix} \begin{bmatrix} \bar{u}_n \\ \hat{u}_n \\ u_n \end{bmatrix}$$

is

$$\begin{aligned} \bar{u}_n &= \bar{u}_0 \sigma^n \\ \hat{u}_n &= [\hat{u}_0 + \bar{u}_0 n \sigma^{-1}] \sigma^n \\ u_n &= [u_0 + \hat{u}_0 n \sigma^{-1} + \bar{u}_0 n(n-1) \sigma^{-2}] \sigma^n \end{aligned} \quad (6.15)$$

6.3 Solution of the Representative O Δ E's

6.3.1 The Operational Form and its Solution

Examples of the nonhomogeneous, linear, first-order ordinary difference equations, produced by applying a time-marching method to the representative equation, are given by Eqs. 6.7 to 6.9. Using the displacement operator, E , these equations can be written

$$[E - (1 + \lambda h)]u_n = h \cdot a e^{\mu h n} \quad (6.16)$$

$$[(1 - \lambda h)E - 1]u_n = h \cdot E \cdot a e^{\mu h n} \quad (6.17)$$

$$\begin{bmatrix} E & -(1 + \lambda h) \\ -\frac{1}{2}(1 + \lambda h)E & E - \frac{1}{2} \end{bmatrix} \begin{bmatrix} \tilde{u} \\ u \end{bmatrix}_n = h \cdot \begin{bmatrix} 1 \\ \frac{1}{2}E \end{bmatrix} \cdot a e^{\mu h n} \quad (6.18)$$

All three of these equations are subsets of the *operational form* of the representative OΔE

$$\boxed{P(E)u_n = Q(E) \cdot ae^{\mu hn}} \quad (6.19)$$

which is produced by applying time-marching methods to the representative ODE, Eq. 4.33. We can express in terms of Eq. 6.19 all manner of standard time-marching methods having multiple time steps and various types of intermediate predictor-corrector families. The terms $P(E)$ and $Q(E)$ are polynomials in E referred to as the *characteristic polynomial* and the *particular polynomial*, respectively.

The general solution of Eq. 6.19 can be expressed as

$$u_n = \sum_{k=1}^K c_k(\sigma_k)^n + ae^{\mu hn} \cdot \frac{Q(e^{\mu h})}{P(e^{\mu h})} \quad (6.20)$$

where σ_k are the K roots of the characteristic polynomial, $P(\sigma) = 0$. When determinants are involved in the construction of $P(E)$ and $Q(E)$, as would be the case for Eq. 6.18, the ratio $Q(E)/P(E)$ can be found by Kramer's rule. Keep in mind that for methods such as in Eq. 6.18 there are multiple (two in this case) solutions, one for u_n and \tilde{u}_n and we are usually only interested in the final solution u_n . Notice also, the important subset of this solution which occurs when $\mu = 0$, representing a time invariant particular solution, or a steady state. In such a case

$$u_n = \sum_{k=1}^K c_k(\sigma_k)^n + a \cdot \frac{Q(1)}{P(1)}$$

6.3.2 Examples of Solutions to Time-Marching OΔE's

As examples of the use of Eqs. 6.19 and 6.20, we derive the solutions of Eqs. 6.16 to 6.18. For the explicit Euler method, Eq. 6.16, we have

$$\begin{aligned} P(E) &= E - 1 - \lambda h \\ Q(E) &= h \end{aligned} \quad (6.21)$$

and the solution of its representative OΔE follows immediately from Eq. 6.20:

$$u_n = c_1(1 + \lambda h)^n + ae^{\mu hn} \cdot \frac{h}{e^{\mu h} - 1 - \lambda h}$$

For the implicit Euler method, Eq. 6.17, we have

$$\begin{aligned} P(E) &= (1 - \lambda h)E - 1 \\ Q(E) &= hE \end{aligned} \quad (6.22)$$

so

$$u_n = c_1 \left(\frac{1}{1 - \lambda h} \right)^n + a e^{\mu h n} \cdot \frac{h e^{\mu h}}{(1 - \lambda h) e^{\mu h} - 1}$$

In the case of the coupled predictor-corrector equations, Eq. 6.18, one solves for the final family u_n (one can also find a solution for the intermediate family \tilde{u}), and there results

$$P(E) = \det \begin{bmatrix} E & -(1 + \lambda h) \\ -\frac{1}{2}(1 + \lambda h)E & E - \frac{1}{2} \end{bmatrix} = E \left(E - 1 - \lambda h - \frac{1}{2} \lambda^2 h^2 \right)$$

$$Q(E) = \det \begin{bmatrix} E & h \\ -\frac{1}{2}(1 + \lambda h)E & \frac{1}{2}hE \end{bmatrix} = \frac{1}{2}hE(E + 1 + \lambda h)$$

The σ -root is found from

$$P(\sigma) = \sigma \left(\sigma - 1 - \lambda h - \frac{1}{2} \lambda^2 h^2 \right) = 0$$

which has only one nontrivial root ($\sigma = 0$ is simply a shift in the reference index). The complete solution can therefore be written

$$u_n = c_1 \left(1 + \lambda h + \frac{1}{2} \lambda^2 h^2 \right)^n + a e^{\mu h n} \cdot \frac{\frac{1}{2}h(e^{\mu h} + 1 + \lambda h)}{e^{\mu h} - 1 - \lambda h - \frac{1}{2} \lambda^2 h^2} \quad (6.23)$$

6.4 The $\lambda - \sigma$ Relation

6.4.1 Establishing the Relation

We have now been introduced to two basic kinds of roots, the λ -roots and the σ -roots. The former are the eigenvalues of the A matrix in the ODE's found by space differencing the original PDE, and the latter are the roots of the characteristic polynomial in a representative ODE found by applying a time-march method to the representative ODE. There is a fundamental relation between the two which can be used to identify many of the essential properties of a time-march method. This relation is first demonstrated by developing it for the explicit Euler method.

First we make use of the semi-discrete approach to find a system of ODE's and then express its solution in the form of Eq. 4.26. Remembering that $t = nh$, one can write

$$u(t) = c_1 (e^{\lambda_1 h})^n \vec{x}_1 + \cdots + c_m (e^{\lambda_m h})^n \vec{x}_m + \cdots + c_M (e^{\lambda_M h})^n \vec{x}_M + P.S. \quad (6.24)$$

where for the present we are not interested in the form of the particular solution. Now the explicit Euler method produces for each λ -root, one σ -root, which is given

by $\sigma = 1 + \lambda h$. So if we use the Euler method for the time advance of the ODE's, the solution¹ of the resulting OΔE is

$$u_n = c_1(\sigma_1)^n \vec{x}_1 + \cdots + c_m(\sigma_m)^n \vec{x}_m + \cdots + c_M(\sigma_M)^n \vec{x}_M + P.S. \quad (6.25)$$

where the c_m and the \vec{x}_m in the two equations are identical and $\sigma_m = (1 + \lambda_m h)$. Comparing Eq. 6.24 and Eq. 6.25, we see a correspondance between σ_m and $e^{\lambda_m h}$. Since the value of $e^{\lambda h}$ can be expressed in terms of the series

$$e^{\lambda h} = 1 + \lambda h + \frac{1}{2}\lambda^2 h^2 + \frac{1}{6}\lambda^3 h^3 + \cdots + \frac{1}{n!}\lambda^n h^n + \cdots$$

the truncated expansion $\sigma = 1 + \lambda h$ is a reasonable² approximation for small enough λh .

Suppose, instead of the Euler method, we use the leapfrog method for the time advance, which is defined by

$$u_{n+1} = u_{n-1} + 2hu'_n \quad (6.26)$$

Applying Eq. 6.6 to Eq. 6.26, we have the characteristic polynomial $P(E) = E^2 - 2\lambda h E - 1$, so that for every λ the σ must satisfy the relation

$$\sigma_m^2 - 2\lambda_m h \sigma_m - 1 = 0 \quad (6.27)$$

Now we notice that each λ produces two σ -roots. For one of these (see the following discussion of spurious roots for the other) we find

$$\sigma_m = \lambda_m h + \sqrt{1 + \lambda_m^2 h^2} \quad (6.28)$$

$$= 1 + \lambda_m h + \frac{1}{2}\lambda_m^2 h^2 - \frac{1}{8}\lambda_m^4 h^4 + \cdots \quad (6.29)$$

On the basis of the isolation theorem, we now place this root, instead of $(1 + \lambda_m h_m)$, into Eq. 6.25 and the result is an $O(h^2)$ method with an error $O(\lambda^3 h^3)$. The following generalizes this concept.

¹From the Isolation Theorem.

²The error is $O(\lambda^2 h^2)$.

6.4.2 The Principal σ -Root

Based on the above we make the follow observation.

Application of the same time-marching method to all of the equations in the coupled A matrix of Eq. 4.6, always produces one σ -root for every λ -root that satisfies the relation

$$\sigma = 1 + \lambda h + \frac{1}{2}\lambda^2 h^2 + \cdots + \frac{1}{k!}\lambda^k h^k + O(h^{k+1}) \quad (6.30)$$

where k is the order of the time-march method.

We refer to the root that has the above property as the *principal* σ -root, and designate it $(\sigma_m)_1$. The above property can be stated regardless of the details of the time-march method, knowing only that its leading error is $O(h^{k+1})$.

6.4.3 Spurious σ -Roots

We saw from Eq. 6.27 that the $\lambda - \sigma$ relation for the leapfrog method produces two σ -roots for each λ . One of these we identified as the principal root which always has the property given in 6.30. The other is referred to as a *spurious* σ -root and designated $(\sigma_m)_2$. In general, the $\lambda - \sigma$ relation produced by a time-march scheme can result in multiple σ -roots all of which, except for the principal one, are spurious. All spurious roots are designated $(\sigma_m)_k$ where $k = 2, 3, \dots$. No matter whether a σ -root is principal or spurious, it is always some algebraic function of the product λh . To express this fact we use the notation $\sigma = \sigma(\lambda h)$.

Spurious roots originate entirely from the numerical approximation of the time-march method and have nothing to do with the ODE being solved. However, generation of spurious roots does not, in itself, make a method inferior. In fact, many very accurate methods in practical use for integrating some forms of ODE have spurious roots.

If a time-march method produces spurious σ -roots, the solution for the ODE in the form shown in Eq. 6.25 must be modified. Following again the message of the Isolation Theorem, we have

$$\begin{aligned} u_n = & c_{11}(\sigma_1)_1^n \vec{x}_1 + \cdots + c_{m1}(\sigma_m)_1^n \vec{x}_m + \cdots + c_{M1}(\sigma_M)_1^n \vec{x}_M + P.S. \\ & + c_{12}(\sigma_1)_2^n \vec{x}_1 + \cdots + c_{m2}(\sigma_m)_2^n \vec{x}_m + \cdots + c_{M2}(\sigma_M)_2^n \vec{x}_M \\ & + c_{13}(\sigma_1)_3^n \vec{x}_1 + \cdots + c_{m3}(\sigma_m)_3^n \vec{x}_m + \cdots + c_{M3}(\sigma_M)_3^n \vec{x}_M \\ & + \text{etc., if there are more spurious roots} \end{aligned} \quad (6.31)$$

It should be mentioned that methods with spurious roots are not self starting. For example, if there is one spurious root to a method, all of the coefficients $(c_m)_2$ in Eq. 6.31 must be initialized by some starting procedure. Presumably (i.e., if one starts the method properly) the spurious coefficients are all initialized with very small magnitudes, and presumably the magnitudes of the spurious roots themselves are all less than one (see Chapter 7). Then the presence of spurious roots does not contaminate the answer. That is, after some finite time (large n) the amplitude of the error associated with the spurious roots are driven to zero.

6.4.4 One-Root Time-Marching Methods

There are a number of time-marching methods that produce only one σ -root for each λ -root. We refer to them as *one-root* methods. They are also called one-step methods. They have the significant advantage of being self-starting which carries with it the very useful property that the time-step interval can be changed at will throughout the marching process. Three one-root methods were analyzed in Section 6.3.2. A popular method having this property, the so-called θ -method, is given by the formula

$$u_{n+1} = u_n + h \left[(1 - \theta)u'_n + \theta u'_{n+1} \right]$$

The θ -method represents the explicit Euler ($\theta = 0$), the trapezoidal ($\theta = \frac{1}{2}$), and the implicit Euler methods ($\theta = 1$), respectively. Its $\lambda - \sigma$ relation is

$$\sigma = \frac{1 + (1 - \theta)\lambda h}{1 - \theta\lambda h}$$

It is instructive to compare the exact solutions to a set of ODE's (with a complete eigensystem) having time-invariant forcing terms, with the exact solution to the ODE's for one-root methods. These are

$$\begin{aligned} \vec{u}(t) &= c_1 \left(e^{\lambda_1 h} \right)^n \vec{x}_1 + \cdots + c_m \left(e^{\lambda_m h} \right)^n \vec{x}_m + \cdots + c_M \left(e^{\lambda_M h} \right)^n \vec{x}_M + A^{-1} \vec{f} \\ \vec{u}_n &= c_1 (\sigma_1)^n \vec{x}_1 + \cdots + c_m (\sigma_m)^n \vec{x}_m + \cdots + c_M (\sigma_M)^n \vec{x}_M + A^{-1} \vec{f} \end{aligned} \quad (6.32)$$

respectively. Notice that when t and $n = 0$, these equations are identical, so that all the constants, vectors, and matrices are identical except the \vec{u} and the terms inside the parentheses on the right hand sides. The only error made by introducing the time march is the error that σ makes in approximating $e^{\lambda h}$.

6.5 Accuracy Measures of Time-Marching Methods

6.5.1 Local and Global Error Measures

There are two broad categories of errors that can be used to evaluate time-marching methods. One is the error made in each time step. This is a *local* error such as that found from a Taylor table analysis, see Section 3.4. It is usually used as the basis for establishing the order of a method. The other is the error determined at the end of a given *event* which has covered a specific interval of time composed of many time steps. This is a *global* error. It is useful for comparing methods, as we shall see in Chapter 8.

It is quite common to judge a time-march method on the basis of results found from a Taylor table. However, a Taylor series analysis is a very limited tool for finding the more subtle properties of a numerical time-marching method. For example, it is of no use in:

- finding spurious roots.
- evaluating numerical stability and separating the errors in phase and amplitude.
- analyzing the particular solution of predictor-corrector combinations.
- finding the global error.

The latter three of these are of concern to us here, and to study them we make use of the material developed in the previous sections of this chapter.

6.5.2 Local Accuracy of the Transient Solution $(er_\lambda, |\sigma|, er_\omega)$

Transient error

The particular choice of an error measure, either local or global, is to some extent arbitrary. However, a necessary condition for the choice should be that the measure can be used consistently for all methods. In the discussion of the λ - σ relation we saw that all time-marching methods produce a principal σ -root for every λ -root that exists in a set of linear ODE's. Therefore, a very natural local error measure for the transient solution is the value of the difference between solutions based on these two roots. We designate this by er_λ and make the following definition

$$er_\lambda \equiv e^{\lambda h} - \sigma_1$$

Both terms are expanded in a Taylor series, and the *error* of the method is the first nonvanishing term. The *order* of the method is the last power of λh matched exactly. This is similar to the error found from a Taylor table.

Amplitude and Phase Error

Suppose a λ eigenvalue is imaginary. Such can indeed be the case when we study the equations governing periodic convection which produces harmonic motion. For such cases it is more meaningful to express the error in terms of amplitude and phase. Let $\lambda = i\omega$ where ω is a real number representing a frequency. Then the numerical method must produce a principal σ -root that is complex and expressible in the form

$$\sigma_1 = \sigma_r + i\sigma_i \approx e^{i\omega h} \quad (6.33)$$

From this it follows that the local error in amplitude is measured by the deviation of $|\sigma_1|$ from unity, where

$$|\sigma_1| = \sqrt{\sigma_r^2 + \sigma_i^2}$$

and the local error in phase can be defined as the first nonvanishing term in a Taylor series expansion of

$$er_\omega \equiv \omega h - \tan^{-1}(\sigma_i/\sigma_r) \quad (6.34)$$

Amplitude and phase errors are important measures of the suitability of time-marching methods used to study oscillating flows.

6.5.3 Local Accuracy of the Particular Solution (er_μ)

The numerical error in the particular solution is found by comparing the particular solution of the ODE with that for the O Δ E. We have found these to be given by

$$P.S._{(ODE)} = ae^{\mu t} \cdot \frac{1}{(\mu - \lambda)}$$

and

$$P.S._{(O\Delta E)} = ae^{\mu t} \cdot \frac{Q(e^{\mu h})}{P(e^{\mu h})}$$

respectively. For a measure of the *local* error in the particular solution we introduce the definition

$$er_\mu \equiv h(\mu - \lambda) \left\{ \frac{P.S._{(O\Delta E)}}{P.S._{(ODE)}} - 1 \right\} \quad (6.35)$$

The multiplication by h converts the error from a global measure to a local one, so that the order of er_λ and er_μ are consistent. The multiplication by $(\mu - \lambda)$, which is arbitrary, is used to create a more uniform normalization. Eq. 6.35 can be written in terms of the characteristic and particular polynomials as

$$\boxed{er_\mu = c_o \cdot \left\{ (\mu - \lambda)Q(e^{\mu h}) - P(e^{\mu h}) \right\}} \quad (6.36)$$

where

$$c_o = \left[\frac{h(\mu - \lambda)}{P(e^{\mu h})} \right]_{h=0}$$

The value of c_o is a method-dependent constant that is often equal to one. The error measure is the first nonvanishing term in the Taylor series expansion. If the forcing function is independent of time, μ is equal to zero, and for this case, many numerical methods generate an er_μ that is also zero.

The algebra involved in finding the order of er_μ can be quite tedious. However, this order is quite important in determining the true order of a time-marching method by the process that has been outlined. An illustration of this is given in the section on Runge-Kutta methods.

6.5.4 Time Accuracy For Nonlinear Applications

In practice, time-marching methods are usually applied to nonlinear ODE's, and it is necessary that the advertised order of accuracy be valid for the nonlinear cases as well as for the linear ones. A necessary condition for this to occur is that the local accuracies of *both the transient and the particular solutions be of the same order*. More precisely, a time-marching method is said to be of order k if

$$er_\lambda = c_1 \cdot (\lambda h)^{k_1+1} \quad (6.37)$$

$$er_\mu = c_2 \cdot (\lambda h)^{k_2+1} \quad (6.38)$$

$$\text{where } k = \text{smallest of } (k_1, k_2) \quad (6.39)$$

6.5.5 Global Accuracy

In contrast to the local error measures which have just been discussed, we can also define global error measures. These are useful when we come to the evaluation of time-marching methods for specific purposes. This subject is covered in Chapter 8 after our introduction to stability in Chapter 7.

Suppose we wish to compute some time-accurate phenomenon over a fixed interval of time using a constant time step. We refer to such a computation as an "event".

Let T be the fixed time of the event and h be the chosen step size. Then the required number of time steps, is N , given by the relation

$$T = Nh$$

Global error in the transient

A natural extension of er_λ to cover the error in an entire event is given by

$$Er_\lambda \equiv e^{\lambda T} - (\sigma_1(\lambda h))^N \quad (6.40)$$

Global error in amplitude and phase

If the event is periodic, we are more concerned with the global error in amplitude and phase. These are given by how much

$$\left(\sqrt{\sigma_r^2 + \sigma_i^2}\right)^N$$

deviates from unity, and the first nonvanishing term in the expansion of

$$Er_\omega \equiv N \left[\omega h - \tan^{-1} \left(\frac{\sigma_i}{\sigma_r} \right) \right] \quad (6.41)$$

$$= \omega T - N \tan^{-1}(\sigma_i/\sigma_r) \quad (6.42)$$

Global error in the particular solution

Finally, the global error in the particular solution follows naturally by comparing the solutions to the ODE and the OΔE. It can be measured by

$$Er_\mu \equiv (\mu - \lambda) \frac{Q(e^{\mu h})}{P(e^{\mu h})} - 1$$

6.6 Linear Multistep Methods

In the previous sections, we have developed the framework of error analysis for time advance methods and have randomly introduced a few methods without addressing motivational, developmental or design issues. In the subsequent sections, we introduce classes of methods along with their associated error analysis. We shall not spend much time on development or design of these methods, since most of them have historic origins from a wide variety of disciplines and applications. The Linear Multistep Methods (LMM's) are probably the most natural extension to time marching of the space differencing schemes introduced in Chapter 3 and can be analyzed for accuracy or designed using the Taylor Table approach of Section 3.4.

6.6.1 The General Formulation

When applied to the nonlinear ODE

$$\frac{du}{dt} = u' = F(u, t)$$

all linear multistep methods can be expressed³ in the general form

$$\sum_{k=1-K}^1 \alpha_k u_{n+k} = h \sum_{k=1-K}^1 \beta_k F_{n+k} \quad (6.43)$$

where the notation for F is defined in Section 6.1. The methods are said to be linear because the α 's and β 's are independent of u and n , and they are said to be K -step because K time-levels of data are required to march the solution one time-step, h .

When Eq. 6.43 is applied to the representative equation 4.30, and the result is expressed in operational form, one finds

$$\left(\sum_{k=1-K}^1 \alpha_k E^k \right) u_n = h \left(\sum_{k=1-K}^1 \beta_k E^k \right) (\lambda u_n + a e^{\mu h n}) \quad (6.44)$$

We recall from Section 6.5.2 that a time-marching method when applied to the representative equation must provide a σ -root, labeled σ_1 , that approximates $e^{\lambda h}$ through the order of the method. The condition referred to as *consistency* simply means that $\sigma \rightarrow 1$ as $h \rightarrow 0$, and it is certainly a necessary condition for the accuracy of any time marching method. We can also agree that, to be of any value in time accuracy, a method should at least be first-order accurate, that is $\sigma \rightarrow (1 + \lambda h)$ as $h \rightarrow 0$. One can show that these conditions are met by any method represented by Eq. 6.43 if

$$\sum_k \alpha_k = 0 \quad \text{and} \quad \sum_k \beta_k = \sum_k (K + k - 1) \alpha_k$$

Since both sides of Eq. 6.43 can be multiplied by an arbitrary constant, these methods are often “normalized” by requiring

$$\sum_k \beta_k = 1$$

Under this condition $c_o = 1$ in Eq. 6.36.

³Note the shift in indexing between these equations and Eq. 6.2.

6.6.2 Simple Examples

There are many special explicit and implicit forms of linear multistep methods. Two well-known families of them, referred to as Adams-Bashforth (*explicit*) and Adams-Moulton (*implicit*), can be designed using the Taylor Table approach of Section 3.4. Three-step methods from these two families can be written in the following general form

$$u_{n+1} = u_n + h(\beta_1 u'_{n+1} + \beta_0 u'_n + \beta_{-1} u'_{n-1} + \beta_{-2} u'_{n-2}) \quad (6.45)$$

A Taylor Table for Eq. 6.45 can be generated as

	u_n	$h \cdot u'$	$h^2 \cdot u''$	$h^3 \cdot u'''$	$h^4 \cdot u''''$
u_{n+1}	1	1	$\frac{1}{2}$	$\frac{1}{6}$	$\frac{1}{24}$
$-u_n$	-1				
$-h\beta_1 u'_{n+1}$		$-\beta_1$	$-\beta_1$	$-\beta_1 \frac{1}{2}$	$-\beta_1 \frac{1}{6}$
$-h\beta_0 u'_n$		$-\beta_0$			
$-h\beta_{-1} u'_{n-1}$		$-\beta_{-1}$	$-\beta_{-1}$	$-\beta_{-1} \frac{1}{2}$	$-\beta_{-1} \frac{1}{6}$
$-h\beta_{-2} u'_{n-2}$		$-(-2)^0 \beta_{-2}$	$-(-2)^1 \beta_{-2}$	$-(-2)^2 \beta_{-2} \frac{1}{2}$	$-(-2)^3 \beta_{-2} \frac{1}{3}$
	<hr/> 0	<hr/> 0	<hr/> 0	<hr/> 0	<hr/> 0

Table 3.1. Taylor table for Adams-Bashforth/Moulton three-step linear multistep methods.

This leads to the linear system

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 0 & -2 & -4 \\ 3 & 0 & 3 & 12 \\ 4 & 0 & -4 & -32 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_0 \\ \beta_{-1} \\ \beta_{-2} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad (6.46)$$

to solve for the β 's, resulting in a variety of methods of up to fourth-order accuracy. With $\beta_1 \neq 0$, the implicit Adams-Moulton family is obtained, while with $\beta_1 = 0$, the explicit Adams-Bashforth family is obtained. In the latter case, the β 's are obtained from

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & -2 & -4 \\ 0 & 3 & 12 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_{-1} \\ \beta_{-2} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad (6.47)$$

and the maximum order attainable with a three-step method is third order.

A list of simple methods, some of which are very common in CFD applications, is given below together with identifying names that are sometimes associated with them. In the following material AB(n) and AM(n) are used as abbreviations for the (n)th order Adams-Bashforth and (n)th order Adams-Moulton methods. One can verify that the Adams type schemes given below satisfy Eqs. 6.46 and 6.47 up to the order of the method.

Explicit Methods

$$\begin{array}{ll}
 u_{n+1} &= u_n + hu'_n & \text{Euler} \\
 u_{n+1} &= u_{n-1} + 2hu'_n & \text{Leapfrog} \\
 u_{n+1} &= u_n + \frac{1}{2}h[3u'_n - u'_{n-1}] & \text{AB2} \\
 u_{n+1} &= u_n + \frac{h}{12}[23u'_n - 16u'_{n-1} + 5u'_{n-2}] & \text{AB3}
 \end{array}$$

Implicit Methods

$$\begin{array}{ll}
 u_{n+1} &= u_n + hu'_{n+1} & \text{Implicit Euler} \\
 u_{n+1} &= u_n + \frac{1}{2}h[u'_n + u'_{n+1}] & \text{Trapezoidal (AM2)} \\
 u_{n+1} &= \frac{1}{3}[4u_n - u_{n-1} + 2hu'_{n+1}] & \text{2nd O Backward} \\
 u_{n+1} &= u_n + \frac{h}{12}[5u'_{n+1} + 8u'_n - u'_{n-1}] & \text{AM3}
 \end{array}$$

6.6.3 Two-Step Linear Multistep Methods

High resolution CFD problems usually require very large data bases to store the spatial information from which the time derivative is calculated. This limits the interest in multistep methods to about two time levels. The most general two-step linear multistep method (i.e., K=2 in Eq. 6.43), that is at least first-order accurate, can be written as

$$(1 + \xi)u_{n+1} = [(1 + 2\xi)u_n - \xi u_{n-1}] + h[\theta u'_{n+1} + (1 - \theta + \varphi)u'_n - \varphi u'_{n-1}] \quad (6.48)$$

where we have shifted our reference frame so that $n + 1$ represents the most advanced time level. Clearly the methods are explicit if $\theta = 0$ and implicit otherwise. A list of methods contained in Eq. 6.48 is given in Table 6.1.

θ	ξ	φ	Method	Order
0	0	0	Euler	1
1	0	0	Implicit Euler	1
1/2	0	0	Trapezoidal or AM2	2
1	1/2	0	2nd Order Backward	2
3/4	0	-1/4	Adams type	2
1/3	-1/2	-1/3	Lees Type	2
1/2	-1/2	-1/2	Two-step trapezoidal	2
5/9	-1/6	-2/9	A-contractive	2
0	-1/2	0	Leapfrog	2
0	0	1/2	AB2	2
0	-5/6	-1/3	Most accurate explicit	3
1/3	-1/6	0	Third-order implicit	3
5/12	0	1/12	AM3	3
1/6	-1/2	-1/6	Milne	4

Table 6.1. Some linear one- and two-step methods, see Eq. 6.48.

One can show after a little algebra⁴ that both er_μ and er_λ are reduced to $O(h^3)$ (i.e., the methods are 2nd-order accurate) if

$$\varphi = \xi - \theta + \frac{1}{2}$$

The class of all 3rd-order methods is determined by imposing the additional constraint

$$\xi = 2\theta - \frac{5}{6}$$

Finally a unique fourth-order method is found by setting $\theta = -\varphi = -\xi/3 = \frac{1}{6}$.

6.7 Predictor-Corrector Methods

There are a wide variety of predictor-corrector schemes created and used for a variety of purposes. Their use in solving ODE's is relatively easy to illustrate and understand. Their use in solving PDE's can be much more subtle and demands concepts⁵ which have no counterpart in the analysis of ODE's.

⁴A Taylor table is ideal for this.

⁵Such as alternating direction, fractional-step, and hybrid methods.

Predictor-corrector methods constructed to time-march linear or nonlinear ODE's are composed of sequences of linear multistep methods, each of which is referred to as a family in the solution process. There may be many families in the sequence, and usually the final family has a higher Taylor-series order of accuracy than the intermediate ones. Their use is motivated by ease of application and increased efficiency, where measures of efficiency are discussed in the next two chapters.

A simple one-predictor, one-corrector example is given by

$$\begin{aligned}\tilde{u}_{n+\alpha} &= u_n + \alpha h F_n \\ u_{n+1} &= u_n + h [\beta \tilde{F}_{n+\alpha} + \gamma F_n]\end{aligned}\tag{6.49}$$

where the parameters α, β and γ are arbitrary parameters to be determined. One can analyze this sequence by applying it to the representative equation and using the operational techniques outlined in Section 6.3. It is easy to show, following the example leading to Eq. 6.23, that

$$P(E) = E^\alpha \cdot [E - 1 - (\gamma + \beta)\lambda h - \alpha\beta\lambda^2 h^2]\tag{6.50}$$

$$Q(E) = E^\alpha \cdot h \cdot [\beta E^\alpha + \gamma + \alpha\beta\lambda h]\tag{6.51}$$

Considering only local accuracy, one is led, by following the discussion in Section 6.5, to the following observations. For the method to be second-order accurate *both* er_λ and er_μ must be $O(h^3)$. For this to hold for er_λ , it is obvious from Eq. 6.50 that

$$\gamma + \beta = 1 \quad ; \quad \alpha\beta = \frac{1}{2}$$

which provides two equations for three unknowns. The situation for er_μ requires some algebra, but it is not difficult to show using Eq. 6.36 that the same conditions also make it $O(h^3)$. One concludes, therefore, that the predictor-corrector sequence

$$\begin{aligned}\tilde{u}_{n+\alpha} &= u_n + \alpha h u'_n \\ u_{n+1} &= u_n + \frac{1}{2}h \left[\left(\frac{1}{\alpha} \right) \tilde{u}'_{n+\alpha} + \left(\frac{2\alpha - 1}{\alpha} \right) u'_n \right]\end{aligned}\tag{6.52}$$

is a second-order accurate method for any α .

A classical predictor-corrector sequence is formed by following an Adams-Bashforth predictor of any order with an Adams-Moulton corrector having an order one higher. The order of the combination is then equal to the order of the corrector. If the order of the corrector is (k) , we refer to these as ABM(k) methods. The Adams-Bashforth-Moulton sequence for $k = 3$ is

$$\begin{aligned}\tilde{u}_{n+1} &= u_n + \frac{1}{2}h [3u'_n - u'_{n-1}] \\ u_{n+1} &= u_n + \frac{h}{12} [5\tilde{u}'_{n+1} + 8u'_n - u'_{n-1}]\end{aligned}\tag{6.53}$$

Some simple, specific, second-order accurate methods are given below. The Gazdag method, which we discuss in Chapter 8, is

$$\begin{aligned}\tilde{u}_{n+1} &= u_n + \frac{1}{2}h[3\tilde{u}'_n - \tilde{u}'_{n-1}] \\ u_{n+1} &= u_n + \frac{1}{2}h[\tilde{u}'_n + \tilde{u}'_{n+1}]\end{aligned}\tag{6.54}$$

The Burstein method, obtained from Eq. 6.52 with $\alpha = 1/2$ is

$$\begin{aligned}\tilde{u}_{n+1/2} &= u_n + \frac{1}{2}hu'_n \\ u_{n+1} &= u_n + h\tilde{u}'_{n+1/2}\end{aligned}\tag{6.55}$$

and, finally, MacCormack's method, also discussed later, is

$$\begin{aligned}\tilde{u}_{n+1} &= u_n + hu'_n \\ u_{n+1} &= \frac{1}{2}[u_n + \tilde{u}_{n+1} + h\tilde{u}'_{n+1}]\end{aligned}\tag{6.56}$$

It should be pointed out that the last two sequences are only “skeletons” of the methods used in application. The real power and convenience of the MacCormack method, for example, does not emerge until it is presented for hyperbolic PDE's.

6.8 Runge-Kutta Methods

There is a special subset of predictor-corrector methods, referred to as (one-step) Runge-Kutta methods, that produce just one σ -root for each λ -root such that $\sigma(\lambda h)$ corresponds to the Taylor series expansion of $e^{\lambda h}$ out through the order of the method and then truncates. Thus for a Runge-Kutta method of order k (up to 4th order)

$$\sigma = 1 + \lambda h + \frac{1}{2}\lambda^2 h^2 + \cdots + \frac{1}{k!}\lambda^k h^k\tag{6.57}$$

It is not particularly difficult to build this property into a method, but, as we pointed out in Section 6.6.4, it is not sufficient to guarantee k 'th order accuracy for the solution of $u' = F(u, t)$ or for the representative equation. To ensure k 'th order accuracy, the method must further satisfy the constraint that

$$er_\mu = O(h^{k+1})\tag{6.58}$$

and this is much more difficult.

The most widely publicized Runge-Kutta process is the one that leads to the fourth-order method. We present it below in some detail. It is usually introduced in the form

$$\begin{aligned} k_1 &= hF(u_0, t_0) \\ k_2 &= hF(u_0 + \beta k_1, t_0 + \alpha h) \\ k_3 &= hF(u_0 + \beta_1 k_1 + \gamma_1 k_2, t_0 + \alpha_1 h) \\ k_4 &= hF(u_0 + \beta_2 k_1 + \gamma_2 k_2 + \delta_2 k_3, t_0 + \alpha_2 h) \end{aligned}$$

followed by

$$u(t_0 + h) - u(t_0) = \mu_1 k_1 + \mu_2 k_2 + \mu_3 k_3 + \mu_4 k_4 \quad (6.59)$$

However, we prefer to present it using predictor-corrector notation. Thus, a scheme entirely equivalent to 6.59 is

$$\begin{aligned} \hat{u}_{n+\alpha} &= u_n + \beta h u'_n \\ \tilde{u}_{n+\alpha_1} &= u_n + \beta_1 h u'_n + \gamma_1 h \hat{u}'_{n+\alpha} \\ \bar{u}_{n+\alpha_2} &= u_n + \beta_2 h u'_n + \gamma_2 h \hat{u}'_{n+\alpha} + \delta_2 h \tilde{u}'_{n+\alpha_1} \\ u_{n+1} &= u_n + \mu_1 h u'_n + \mu_2 h \hat{u}'_{n+\alpha} + \mu_3 h \tilde{u}'_{n+\alpha_1} + \mu_4 h \bar{u}'_{n+\alpha_2} \end{aligned} \quad (6.60)$$

Appearing in Eqs. 6.59 and 6.60 are a total of 13 parameters which are to be determined such that the method is fourth-order according to the requirements in Eqs. 6.57 and 6.58. First of all, the choices for the time samplings, α , α_1 , and α_2 , are not arbitrary. They must satisfy the relations

$$\begin{aligned} \alpha &= \beta \\ \alpha_1 &= \beta_1 + \gamma_1 \\ \alpha_2 &= \beta_2 + \gamma_2 + \delta_2 \end{aligned} \quad (6.61)$$

The algebra involved in finding algebraic equations for the remaining 10 parameters is not trivial, but the equations follow directly from finding $P(E)$ and $Q(E)$ and then satisfying the conditions in Eqs. 6.57 and 6.58. Using Eq. 6.61 to eliminate the β 's we find from Eq. 6.57 the four conditions

$$\begin{aligned} \mu_1 + \mu_2 + \mu_3 + \mu_4 &= 1 & (1) \\ \mu_2 \alpha + \mu_3 \alpha_1 + \mu_4 \alpha_2 &= 1/2 & (2) \\ \mu_3 \alpha \gamma_1 + \mu_4 (\alpha \gamma_2 + \alpha_1 \delta_2) &= 1/6 & (3) \\ \mu_4 \alpha \gamma_1 \delta_2 &= 1/24 & (4) \end{aligned} \quad (6.62)$$

These four relations guarantee that the five terms in σ exactly match the first 5 terms in the expansion of $e^{\lambda h}$. To satisfy the condition that $er_\mu = O(k^5)$, we have to fulfill four more conditions

$$\begin{aligned}
 \mu_2\alpha^2 + \mu_3\alpha_1^2 + \mu_4\alpha_2^2 &= 1/3 & (3) \\
 \mu_2\alpha^3 + \mu_3\alpha_1^3 + \mu_4\alpha_2^3 &= 1/4 & (4) \\
 \mu_3\alpha^2\gamma_1 + \mu_4(\alpha^2\gamma_2 + \alpha_1^2\delta_2) &= 1/12 & (4) \\
 \mu_3\alpha\alpha_1\gamma_1 + \mu_4\alpha_2(\alpha\gamma_2 + \alpha_1\delta_2) &= 1/8 & (4)
 \end{aligned} \tag{6.63}$$

The number in parentheses at the end of each equation indicates the order that is the basis for the equation. Thus if the first 3 equations in 6.62 and the first equation in 6.63 are all satisfied, the resulting method would be third-order accurate. Note that only the first seven constraints must be satisfied for linear inhomogeneous ODE's. Therefore, the method of analysis based on $P(E)$ and $Q(E)$ will not lead to the eighth constraint.

There are eight equations in 6.62 and 6.63 which must be satisfied by the 10 unknowns. Since the equations are overdetermined, two parameters can be set arbitrarily. Several choices for the parameters have been proposed, but the most popular one is due to Runge. It results in the “standard” fourth-order Runge-Kutta method expressed in predictor-corrector form as

$$\begin{aligned}
 \hat{u}_{n+1/2} &= u_n + \frac{1}{2}hu'_n \\
 \tilde{u}_{n+1/2} &= u_n + \frac{1}{2}h\hat{u}'_{n+1/2} \\
 \bar{u}_{n+1} &= u_n + h\tilde{u}'_{n+1/2} \\
 u_{n+1} &= u_n + \frac{1}{6}h\left[u'_n + 2(\hat{u}'_{n+1/2} + \tilde{u}'_{n+1/2}) + \bar{u}'_{n+1}\right]
 \end{aligned} \tag{6.64}$$

Notice that this represents the simple sequence of conventional linear multistep methods referred to, respectively, as

$$\left. \begin{array}{l} \text{Euler Predictor} \\ \text{Euler Corrector} \\ \text{Leapfrog Predictor} \\ \text{Milne Corrector} \end{array} \right\} \equiv RK4$$

One can easily show that both the Burstein and the MacCormack methods given by Eqs. 6.55 and 6.56 are second-order Runge-Kutta methods, and third-order methods can be derived from Eqs. 6.60 by setting $\mu_4 = 0$ and satisfying only Eqs. 6.62 and the first equation in 6.63. It is clear that for orders one through four, RK methods of order k require k evaluations of the derivative function to advance the solution

one time step. We shall discuss the consequences of this in Chapter 8. Higher-order Runge-Kutta methods can be developed, but they require more derivative evaluations than their order. For example, a fifth-order method requires 6 evaluations to advance the solution one step. In any event, storage requirements probably eliminate the usefulness of the higher-order methods for CFD applications.

6.9 Implementation of Implicit Methods

We have presented a wide variety of time-marching methods and shown how to derive their $\lambda - \sigma$ relations. In the next chapter, we will see that these methods can have widely different properties with respect to stability. This leads to various trade-offs which must be considered in selecting a method for a specific application. Our presentation of the time-marching methods in the context of a linear scalar equation obscures some of the issues involved in implementing an implicit method for systems of equations and nonlinear equations. These are covered in this Section.

6.9.1 Application to Systems of Equations

Consider first the numerical solution of our representative ODE

$$\frac{du}{dt} = \lambda u + ae^{\mu t} \quad (6.65)$$

using the implicit Euler method. Following the steps outlined in Section 6.1, we obtained

$$(1 - \lambda h)u_{n+1} - u_n = he^{\mu h} \cdot ae^{\mu hn} \quad (6.66)$$

Solving for u_{n+1} gives

$$u_{n+1} = \frac{1}{1 - \lambda h}(u_n + he^{\mu h} \cdot ae^{\mu hn}) \quad (6.67)$$

This calculation does not seem particularly onerous in comparison with the application of an explicit method to this ODE, requiring only an additional division.

Now let us apply the implicit Euler method to our generic system of equations given by

$$\frac{du}{dt} = Au - f(t) \quad (6.68)$$

where u and f are vectors and we still assume that A is not a function of u or t . Now the equivalent to Eq. 6.66 is

$$(1 - hA)u_{n+1} - u_n = hf(t + h) \quad (6.69)$$

or

$$u_{n+1} = (1 - hA)^{-1}[u_n + hf(t + h)] \quad (6.70)$$

Therefore, the simple division required in the scalar case has been replaced by the solution of a linear system of equations. For our one-dimensional examples, the system of equations which must be solved is tridiagonal, and hence its solution is inexpensive, but in multidimensions the bandwidth can be very large. In general, the cost per time step of an implicit method is much larger than that of an explicit method. The primary area of application of implicit methods is in the solution of *stiff* ODE's, as we shall see in Chapter 8.

6.9.2 Application to Nonlinear Equations

Now consider the general *nonlinear* scalar ODE given by

$$\frac{du}{dt} = F(u, t) \quad (6.71)$$

Application of the implicit Euler method gives

$$u_{n+1} = u_n + F(u_{n+1}, t_{n+1}) \quad (6.72)$$

This is a nonlinear $O\Delta E$. As an example, consider the nonlinear ODE

$$\frac{du}{dt} + \frac{1}{2}u^2 = 0 \quad (6.73)$$

solved using implicit Euler time marching which gives

$$u_{n+1} + h\frac{1}{2}u_{n+1}^2 = u_n \quad (6.74)$$

which requires a nontrivial method to solve for u_{n+1} . However, in practice, the “initial guess” for this nonlinear problem is quite close to the solution, since the “initial guess” is simply the solution at the previous time step, which implies that a linearization approach may be quite successful. There are several different approaches one can take to solving this nonlinear difference equation. An iterative method, such as Newton's method (see below), can be used. Alternatively, the nonlinear ODE, Eq. 6.71, can be linearized as described in the next Section.

6.9.3 Local Linearization for Scalar Equations

General Development

Let us start the process of local linearization by considering Eq. 6.71. In order to implement the linearization, we expand $F(u, t)$ about some reference point in time. Designate the reference value by t_n and the corresponding value of the dependent variable by u_n . A Taylor series expansion about these reference quantities gives

$$\begin{aligned}
 F(u, t) = & F(u_n, t_n) + \left(\frac{\partial F}{\partial u} \right)_n (u - u_n) + \left(\frac{\partial F}{\partial t} \right)_n (t - t_n) \\
 & + \frac{1}{2} \left(\frac{\partial^2 F}{\partial u^2} \right)_n (u - u_n)^2 + \left(\frac{\partial^2 F}{\partial u \partial t} \right)_n (u - u_n)(t - t_n) \\
 & + \frac{1}{2} \left(\frac{\partial^2 F}{\partial t^2} \right)_n (t - t_n)^2 + \dots
 \end{aligned} \tag{6.75}$$

On the other hand, the expansion of $u(t)$ in terms of the independent variable t is

$$u(t) = u_n + (t - t_n) \left(\frac{\partial u}{\partial t} \right)_n + \frac{1}{2} (t - t_n)^2 \left(\frac{\partial^2 u}{\partial t^2} \right)_n + \dots \tag{6.76}$$

If t is within h of t_n , both $(t - t_n)^k$ and $(u - u_n)^k$ are $O(h^k)$, and Eq. 6.75 can be written

$$F(u, t) = F_n + \left(\frac{\partial F}{\partial u} \right)_n (u - u_n) + \left(\frac{\partial F}{\partial t} \right)_n (t - t_n) + O(h^2) \tag{6.77}$$

Notice that this is an expansion of the *derivative* of the function. Thus, relative to the order of expansion of the function, it represents a second-order-accurate, locally-linear approximation to $F(u, t)$ that is valid in the vicinity of the reference station t_n and the corresponding $u_n = u(t_n)$. With this we obtain the locally (in the neighborhood of t_n) time-linear representation of Eq. 6.71, namely

$$\frac{du}{dt} = \left(\frac{\partial F}{\partial u} \right)_n u + \left(F_n - \left(\frac{\partial F}{\partial u} \right)_n u_n \right) + \left(\frac{\partial F}{\partial t} \right)_n (t - t_n) + O(h^2) \tag{6.78}$$

Implementation of the Trapezoidal Method

As an example of how such an expansion can be used, consider the mechanics of applying the trapezoidal method for the time march of Eq. 6.71. The trapezoidal method is given by

$$u_{n+1} = u_n + \frac{1}{2}h[F_{n+1} + F_n] + hO(h^2) \quad (6.79)$$

where we write $hO(h^2)$ to emphasize that the method is second order accurate. Using Eq. 6.77 to evaluate $F_{n+1} = F(u_{n+1}, t_{n+1})$, one finds

$$u_{n+1} = u_n + \frac{1}{2}h \left[F_n + \left(\frac{\partial F}{\partial u} \right)_n (u_{n+1} - u_n) + h \left(\frac{\partial F}{\partial t} \right)_n + O(h^2) + F_n \right] + hO(h^2) \quad (6.80)$$

Note that the $O(h^2)$ term within the brackets (which is due to the local linearization) is multiplied by h and therefore is the same order as the $hO(h^2)$ error from the Trapezoidal Method. The use of local time linearization updated at the end of each time step, and the trapezoidal time march, combine to make a *second-order-accurate* numerical integration process. There are, of course, other second-order implicit time-marching methods that can be used. The important point to be made here is that *local linearization updated at each time step has not reduced the order of accuracy* of a second-order time-marching process.

A very useful reordering of the terms in Eq. 6.80 results in the expression

$$\left[1 - \frac{1}{2}h \left(\frac{\partial F}{\partial u} \right)_n \right] \Delta u_n = hF_n + \frac{1}{2}h^2 \left(\frac{\partial F}{\partial t} \right)_n \quad ; \quad O(h^2) \quad (6.81)$$

which is now in the delta form which will be formally introduced in Section 12.6. In many fluid mechanic applications the nonlinear function F is not an *explicit* function of t . In such cases the partial derivative of $F(u)$ with respect to t is zero and Eq. 6.81 simplifies to the second-order accurate expression

$$\left[1 - \frac{1}{2}h \left(\frac{\partial F}{\partial u} \right)_n \right] \Delta u_n = hF_n \quad ; \quad O(h^2) \quad (6.82)$$

Notice that the RHS is extremely simple. It is the product of h and the RHS of the basic equation evaluated at the previous time step. In this example the basic equation was the simple scalar equation 6.71, but for our applications, it is generally the space-differenced form of the steady-state equation of some fluid flow problem.

A numerical time-march procedure using Eq. 6.82 is usually implemented as follows:

1. Solve for the elements of hF_n , store them in an array say R , and save u_n .
2. Solve for the elements of the matrix multiplying Δu_n and store in some appropriate manner making use of sparseness or bandedness of the matrix if possible. Let this storage area be referred to as B .
3. Solve the coupled set of linear equations

$$B\Delta u_n = R$$

for Δu_n . (Very seldom does one find B^{-1} in carrying out this step).

4. Find u_{n+1} by adding Δu_n to u_n , thus

$$u_{n+1} = \Delta u_n + u_n$$

The solution for u_{n+1} probably is stored such that it overwrites the value of u_n and the process is repeated.

Implementation of the Implicit Euler Method

We have seen that the 1'st order Euler implicit method can be written

$$u_{n+1} = u_n + hF_{n+1} \quad ; \quad O(h) \quad (6.83)$$

if we introduce Eq. 6.78 into this method, rearrange terms, and remove the explicit dependency on time, we arrive at the form

$$\left[1 - h \left(\frac{\partial F}{\partial u} \right)_n \right] \Delta u_n = hF_n \quad ; \quad O(h) \quad (6.84)$$

We see that the only difference between the implementation of the trapezoidal method and the implicit Euler method is the factor of $\frac{1}{2}$ in the brackets of the left side of Eqs. 6.82 and 6.84. Omission of this factor degrades the method in time accuracy by one order of h .

Newton's Method

Consider the limit $h \rightarrow \infty$ of Eq. 6.84 obtained by dividing both sides by h and setting $1/h = 0$. There results

$$-\left(\frac{\partial F}{\partial u}\right)_n \Delta u_n = F_n \quad (6.85)$$

or

$$u_{n+1} = u_n - \left[\left(\frac{\partial F}{\partial u} \right)_n \right]^{-1} F_n \quad (6.86)$$

This is the well-known Newton method for finding the roots of a nonlinear equation $F(u) = 0$. The fact that it has quadratic convergence is verified by a glance at Eqs. 6.75 and 6.76 (remember the dependence on t has been eliminated for this case). By quadratic convergence, we mean that the error after a given iteration is proportional to the square of the error at the previous iteration, where the error is the difference between the current solution and the converged solution. Quadratic convergence is thus a very powerful property. Use of a finite value of h in Eq. 6.84 leads to linear convergence, i.e., the error at a given iteration is some multiple of the error at the previous iteration. The reader should ponder the meaning of letting $h \rightarrow \infty$ for the trapezoidal method, given by Eq. 6.82.

6.9.4 Local Linearization for Coupled Sets of Nonlinear Equations

In order to present this concept, let us consider an example involving some simple boundary-layer equations. We choose the Falkner-Skan equations from classical boundary-layer theory. Our task is to apply the implicit trapezoidal method to the equations

$$\frac{d^3 f}{dt^3} + f \frac{d^2 f}{dt^2} + \beta \left(1 - \left(\frac{df}{dt} \right)^2 \right) = 0 \quad (6.87)$$

Here f represents a dimensionless stream function and β is a scaling factor.

First of all we reduce Eq. 6.87 to a set of first-order nonlinear equations by the transformations

$$u_1 = \frac{d^2 f}{dt^2} \quad , \quad u_2 = \frac{df}{dt} \quad , \quad u_3 = f \quad (6.88)$$

This gives the coupled set of three nonlinear equations

$$\begin{aligned}
u'_1 &= F_1 = -u_1 u_3 - \beta (1 - u_2^2) \\
u'_2 &= F_2 = u_1 \\
u'_3 &= F_3 = u_2
\end{aligned} \tag{6.89}$$

and these can be represented in vector notation as

$$\frac{d\vec{u}}{dt} = \vec{F}(\vec{u}) \tag{6.90}$$

Now we seek to make the same local expansion that derived Eq. 6.78, except that this time we are faced with a nonlinear vector function, rather than a simple nonlinear scalar function. The required extension requires the evaluation of a matrix, called the Jacobian matrix. Let us refer to this matrix as A . It is derived from Eq. 6.90 by the following process

$$A = (a_{ij}) = \partial F_i / \partial u_j \tag{6.91}$$

For the general case involving a third order matrix this is

$$A = \begin{bmatrix} \frac{\partial F_1}{\partial u_1} & \frac{\partial F_1}{\partial u_2} & \frac{\partial F_1}{\partial u_3} \\ \frac{\partial F_2}{\partial u_1} & \frac{\partial F_2}{\partial u_2} & \frac{\partial F_2}{\partial u_3} \\ \frac{\partial F_3}{\partial u_1} & \frac{\partial F_3}{\partial u_2} & \frac{\partial F_3}{\partial u_3} \end{bmatrix} \tag{6.92}$$

The expansion of $\vec{F}(\vec{u})$ about some reference state \vec{u}_n can be expressed in a way similar to the scalar expansion given by eq 6.75. Omitting the explicit dependency on the independent variable t , and defining \vec{F}_n as $\vec{F}(\vec{u}_n)$, one has ⁶

$$\vec{F}(\vec{u}) = \vec{F}_n + A_n(\vec{u} - \vec{u}_n) + O(h^2) \tag{6.93}$$

where $t - t_n$ and the argument for $O(h^2)$ is the same as in the derivation of Eq. 6.76. Using this we can write the local linearization of Eq. 6.90 as

$$\frac{d\vec{u}}{dt} = A_n \vec{u} + \underbrace{\left[\vec{F}_n - A_n \vec{u}_n \right]}_{\text{"constant"}} + O(h^2) \tag{6.94}$$

⁶The Taylor series expansion of a vector contains a vector for the first term, a matrix times a vector for the second term, and tensor products for the terms of higher order.

which is a locally-linear, second-order-accurate approximation to a set of coupled nonlinear ordinary differential equations that is valid for $t \leq t_n + h$. Any first- or second-order time-marching method, explicit or implicit, could be used to integrate the equations without loss in accuracy with respect to order. The number of times, and the manner in which, the terms in the Jacobian matrix are updated as the solution proceeds depends, of course, on the nature of the problem.

Returning to our simple boundary-layer example, which is given by Eq. 6.89, we find the Jacobian matrix to be

$$A = \begin{bmatrix} -u_3 & 2\beta u_2 & -u_1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad (6.95)$$

The student should be able to derive results for this example that are equivalent to those given for the scalar case in Eq. 6.81. Thus for the Falkner-Skan equations the trapezoidal method results in

$$\begin{bmatrix} 1 + \frac{h}{2}(u_3)_n & -\beta h(u_2)_n & \frac{h}{2}(u_1)_n \\ -\frac{h}{2} & 1 & 0 \\ 0 & -\frac{h}{2} & 1 \end{bmatrix} \begin{bmatrix} (\Delta u_1)_n \\ (\Delta u_2)_n \\ (\Delta u_3)_n \end{bmatrix} = h \begin{bmatrix} -(u_1 u_3)_n - \beta(1 - u_2^2)_n \\ (u_1)_n \\ (u_2)_n \end{bmatrix} \quad (6.96)$$

We find \vec{u}_{n+1} from $\Delta \vec{u}_n + \vec{u}_n$, and the solution is now advanced one step. Re-evaluate the elements using \vec{u}_{n+1} and continue. Without any iterating within a step advance, the solution will be second-order-accurate in time.